# eScholarship@UMassChan

## Chromosome-level assembly of the Atlantic silverside genome reveals extreme levels of sequence diversity and structural genetic variation [preprint]

| | |
|---|---|
| Item Type | Preprint |
| Authors | Tigano, Anna;Jacobs, Arne;Wilder, Aryn P.;Nand, Ankita;Zhan, Ye;Dekker, Job |
| Citation | <p>bioRxiv 2020.10.27.357293; doi: https://doi.org/10.1101/2020.10.27.357293. <a href="https://doi.org/10.1101/2020.10.27.357293" target="_blank" title="preprint in bioRxiv">Link to preprint on bioRxiv</a>.</p> |
| DOI | 10.1101/2020.10.27.357293 |
| Rights | The copyright holder for this preprint is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY-NC-ND 4.0 International license. |
| Download date | 2024-12-26 00:43:48 |
| Item License | http://creativecommons.org/licenses/by-nc-nd/4.0/ |
| Link to Item | https://hdl.handle.net/20.500.14038/29631 |

1    **Chromosome-level assembly of the Atlantic silverside genome reveals extreme levels of**

2    **sequence diversity and structural genetic variation**

3

4    Anna Tigano[1,2], Arne Jacobs[1], Aryn P. Wilder[1,3], Ankita Nand[4], Ye Zhan[4], Job Dekker[4,5], Nina

5    O. Therkildsen[1]

6    [1]Department of Natural Resources, Cornell University, Ithaca, NY, USA

7    [2]Department of Molecular, Cellular and Biomedical Sciences, University of New Hampshire,

8    Durham, NH, USA

9    [3]Conservation Genetics, San Diego Zoo Global, Escondido, CA, USA

10   [4] Program in Systems Biology, University of Massachusetts Medical School, Worcester, MA

11   01605, USA

12   [5] Howard Hughes Medical Institute, Chevy Chase, MD 20815, USA

13

14

15 **Abstract**

16 The levels and distribution of standing genetic variation in a genome can provide a wealth of

17 insights about the adaptive potential, demographic history, and genome structure of a population

18 or species. As structural variants are increasingly associated with traits important for adaptation

19 and speciation, investigating both sequence and structural variation is essential for wholly

20 tapping this potential. Using a combination of shotgun sequencing, 10X Genomics linked reads

21 and proximity-ligation data (Chicago and Hi-C), we produced and annotated a chromosome-level

22 genome assembly for the Atlantic silverside (*Menidia menidia*) - an established ecological model

23 for studying the phenotypic effects of natural and artificial selection - and examined patterns of

24 genomic variation across two individuals sampled from different populations with divergent

25 local adaptations. Levels of diversity varied substantially across each chromosome, consistently

26 being highly elevated near the ends (presumably near telomeric regions) and dipping to near zero

27 around putative centromeres. Overall, our estimate of the genome-wide average heterozygosity

28 in the Atlantic silverside is the highest reported for a fish, or any vertebrate, to date (1.32-1.76%

29 depending on inference method and sample). Furthermore, we also found extreme levels of

30 structural variation, affecting ~23% of the total genome sequence, including multiple large

31 inversions (> 1 Mb and up to 12.6 Mb) associated with previously identified haploblocks

32 showing strong differentiation between locally adapted populations. These extreme levels of

33 standing genetic variation are likely associated with large effective population sizes and may

34 help explain the remarkable adaptive divergence among populations of the Atlantic silverside.

35

36

37

2

38  **Introduction**

39  Standing genetic variation is widely recognized as the main source of adaptation (Barrett &

40  Schluter 2008; Tigano & Friesen 2016) and is important for natural populations to maximize

41  their potential to adapt to changes in their environment. As genetic diversity is the result of the

42  interplay of mutation, selection, drift and gene flow, the levels and patterns of standing genetic

43  variation found within a species can provide important insights not only about its adaptive

44  potential but also about its demographic and evolutionary history.

45      Traditionally, quantification of standing genetic variation has been based on sequence

46  variation, often across a limited number of genetic markers, or small microsatellite repeats. As an

47  increasing number of empirical studies shows the mosaic nature of the genome (Pääbo 2003)

48  with different genomic regions showing vastly different levels of diversity and differentiation

49  (e.g., Martinez Barrio et al. 2016; Campagna et al. 2017; Murray et al. 2017; Sardell et al. 2018),

50  it is evident that small marker panels do not grant the resolution to describe variation in diversity

51  across the genome (Dutoit et al. 2016). Furthermore, structural variation, including changes in

52  the position, orientation, and number of copies of DNA sequence, is generally neglected as a

53  type of standing genetic variation. Structural variation has been associated directly or indirectly

54  with many traits involved in speciation and adaptation and is abundant in the few genomes in

55  which they have been catalogued (Wellenreuther & Bernatchez 2018; Catanach et al. 2019;

56  Lucek et al. 2019; Mérot et al. 2020; Tigano et al. 2020; Weissensteiner et al. 2020). Structural

57  variants can directly affect phenotypic traits, such as the insertion of a repeated transposable

58  element in the iconic case of industrial melanism in the peppered moth (*Biston betularia*; Van't

59  Hof et al. 2016), or may promote the maintenance of divergent haplotypes between locally

60  adapted populations or groups (e.g. ecotypes or morphs) within single populations via

3

61    recombination suppression (e.g., Faria et al. 2019; Kess et al. 2020). Structural variation is

62    therefore a key source of standing genetic variation, which can also play an important role in

63    rapid evolutionary responses to environmental change (Reid et al. 2016). To better assess levels

64    of standing variation and understand how demographic and evolutionary factors contribute to

65    their distribution in the genome, we need to examine large proportions of the genome, preferably

66    its entirety, and examine sequence and structural variation jointly. A high-quality reference

67    genome for the species of interest is therefore fundamental as we need both broad coverage to

68    accurately assess variation in levels of standing sequence variation across the genome, and high

69    contiguity to investigate standing structural variation.

70        The Atlantic silverside (*Menidia menidia*), a small coastal fish distributed along the

71    Atlantic coast of North America, shows a remarkable degree of local adaptation in a suite of

72    traits, including growth rate, number of vertebrae, and temperature-dependent sex determination

73    (Hice et al. 2012), that are associated with strong environmental gradients across its wide

74    latitudinal range. This species also provided the first discovery of temperature-dependent sex

75    determination in fishes (Conover & Kynard 1981) and was one of the first species in which

76    countergradient phenotypic variation was documented (Conover & Present 1990). Through

77    extensive prior work, the Atlantic silverside has, in fact, become an important ecological model

78    to study the phenotypic effects of selection, both natural and artificial, in the wild and under

79    controlled conditions in the lab (Conover & Munch 2002; Conover et al. 2005; Hice et al. 2012).

80    In one iconic experiment, wild-caught Atlantic silversides were subjected to different size-

81    selective regimes to investigate the potential of fisheries to induce evolutionary change in

82    harvested species (Conover & Munch 2002). Seventeen years later, genomic analysis of fish

83    from that experiment identified substantial allele frequency shifts associated with rapid

4

84    phenotypic shifts in growth rates (Therkildsen et al. 2019). In the absence of a reference genome,

85    genomic reads were mapped to the silverside reference transcriptome, so only protein-coding

86    regions of the genome were analyzed ('in-silico' exome capture). Yet, anchoring the

87    transcriptome contigs to the medaka (*Oryzias latipes*) chromosome-level reference genome

88    revealed that the most conspicuous allele frequency shifts clustered into a single block on

89    chromosome 24, where more than 9,000 SNPs in strong linkage disequilibrium (LD) increased

90    from low (< 0.05) to high frequency (~0.6) in only five generations. Additional data from natural

91    populations across the geographical distribution of the species showed that this same block,

92    likely spanning several Mb of the chromosome, was fixed for opposite haplotypes among wild

93    silverside populations that naturally differ in growth rates (Conover & Present 1990; Conover &

94    Munch 2002; Therkildsen et al. 2019). Moreover, three additional blocks comprising hundreds of

95    genes in high linkage disequilibrium (LD) were found to be segregating among the natural

96    populations, with each LD block ('haploblocks' hereafter) mapping predominantly to unique

97    medaka chromosomes (Wilder et al. 2020). Similar to the haploblock on chromosome 24,

98    opposite haplotypes in these haploblocks were nearly fixed between natural populations that

99    otherwise showed low genome-wide pairwise differentiation. Furthermore, strong LD between

100    genes in these blocks suggested that local recombination suppression, possibly due to inversions,

101    and natural selection maintained these divergent haploblocks in the face of gene flow. It thus

102    appears that large haploblocks play an important role in maintaining local adaptations in the

103    Atlantic silverside, although the exact extent of the genome spanned by these haploblocks and

104    the genomic mechanism maintaining LD are unknown.

105          Given the wealth of ecological information available for the Atlantic silverside and its

106    potential as an evolutionary model to study adaptation and fishery-induced evolutionary change,

107    developing genomic resources for this species is timely and holds great potential for addressing

108    many pressing questions in evolutionary and conservation biology. Previous population genomic

109    analyses based on the transcriptome reference anchored to the medaka genome were limited to

110    the coding genes and, given the unknown degree of synteny conservation between the Atlantic

111    silverside and the medaka, how variants relevant to adaptation and fishery-induced selection

112    clustered in the genome was uncertain. To enable analysis of both coding and non-coding

113    regions, to accurately estimate levels and the genomic distribution of standing genetic variation,

114    both sequence and structural, and to reconstruct the specific genomic structure of the Atlantic

115    silverside genome, we produced a chromosome-level genome assembly for the species using a

116    combination of genomic approaches. Because of known geographic differentiation, we estimated

117    levels of sequence variation within genomes from both the southern and northern parts of the

118    distribution and characterized standing structural variation between these two genomes. Finally,

119    we tested whether the haploblocks identified on four different chromosomes between southern

120    and northern populations were associated with large inversions as the patterns of differentiation

121    and LD suggested (Therkildsen et al. 2019). Our work illustrates the wealth of information that

122    can be obtained from the analysis of one or two genomes in the presence of a high quality

123    reference sequence, and shows that, to the best of our knowledge, the Atlantic silverside has the

124    highest nucleotide diversity reported for a vertebrate to date, and extreme levels of structural

125    variation between two locally adapted populations. The distribution of diversity across the

126    genome is strongly affected by structural variants and, seemingly, by genome features such as

127    centromeres and telomeres. These results taken together highlight the importance of high-quality

128    genomic resources as they enable the joint analysis of sequence and structural variation at the

129    whole-genome level.

6

130 **Methods**

131 *Reference genome assembly*

132 We built a reference genome for the Atlantic silverside through three steps: First, we created a

133 draft assembly using 10X Genomics linked-reads technology (10X Genomics, Pleasanton, CA,

134 USA); second, we used proximity ligation data - Chicago® (Putnam et al. 2016) and Dovetail™

135 Hi-C (Lieberman-Aiden et al. 2009) - from Dovetail Genomics to increase contiguity, break up

136 mis-joins, and orient and join scaffolds into chromosomes; and finally, we used short-insert reads

137 to close gaps in the scaffolded and error-corrected assembly. The data were generated from

138 muscle tissue dissected from two lab-reared F1 offspring of Atlantic silversides collected from

139 the wild on Jekyll Island, Georgia, USA (N 31.02, W 81.43; the southern end of the species

140 distribution range) in May 2017. For 10X Genomics library preparation, we extracted DNA from

141 fresh tissue from one individual using the MagAttract HMW DNA Kit (Qiagen). Prior to library

142 preparation, we selected fragments longer than 30 kb using a BluePippin device (Sage Science).

143 A 10X Genomics library was prepared following standard procedure and sequenced using two

144 lanes of paired-end 150 bp reads on a HiSeq2500 (rapid run mode) at the Biotechnology

145 Resource Center Genomics Facility at Cornell University. To assemble the linked reads, we ran

146 the program *Supernova* (Weisenfeld et al. 2017) from 10X Genomics with varying numbers of

147 reads and compared assembly statistics to identify the number of reads that resulted in the most

148 contiguous assembly. Tissue from the second individual was flash-frozen in liquid nitrogen and

149 shipped to Dovetail Genomics, where Chicago and Hi-C libraries were prepared for further

150 scaffolding. These long-range libraries were sequenced in one lane of Illumina HiSeq X using

151 paired-end 150 bp reads. Two rounds of scaffolding with *HiRise™*, a software pipeline

152 developed specifically for genome scaffolding with Chicago and Hi-C data, were run to scaffold

7

153    and error-correct the best 10X Genomics draft assembly using Dovetail long-range data. Finally,

154    the barcode-trimmed 10X Genomics reads were used to close gaps between contigs.

155        For each of the intermediate and the final assemblies we produced genome contiguity and

156    other assembly statistics using the *assemblathon_stats.pl* script from the Korf Laboratory

157    (https://github.com/KorfLab/Assemblathon/blob/master/assemblath on_stats.pl) and assessed

158    assembly completeness with *BUSCO v3* (Simão et al. 2015) using the Actinopterygii gene set

159    (4584 genes).

160        We estimated the genome size and heterozygosity (i.e. the nucleotide diversity $\pi$ within a

161    single individual) from the raw 10X Genomics data using a k-mer distribution approach. We

162    removed barcodes with the program *longranger basic*, trimmed all reads to the same length of

163    128 bp (as read length is in the equation to estimate genome size) with *cutadapt* (Martin 2011),

164    and estimated the distribution of 25-mers using *Jellyfish* (Marçais & Kingsford 2011). Finally,

165    we analyzed the 25-mers distribution with the web application of *GenomeScope* (Vurture et al.

166    2017), which runs mixture models based on the binomial distributions of k-mer profiles to

167    estimate genome size, heterozygosity and repeat content.

168

169    *Synteny with medaka*

170    The chromosome-level genome assembly of medaka (*Oryzias latipes*) was used by Therkildsen

171    et al. (2019) to order and orient contigs of the Atlantic silverside transcriptome (Therkildsen &

172    Baumann 2020). Although the two species carry the same number of chromosomes (Uwa &

173    Ojima 1981; Warkentine et al. 1987) and few interchromosomal rearrangements have been

174    observed between other species within the Atherinomorpha clade (Amores et al. 2014; Miller et

175    al. 2019), the estimated divergence time between medaka and Atlantic silverside is 91 million

176    years (estimate based on 15 studies, timetree.org) and the degree of syntenic conservation

177    between the two species was unknown. We assessed synteny between the two species using the

178    newly assembled Atlantic silverside reference genome. We aligned the silverside genome to the

179    medaka genome (GenBank assembly accession GCA_002234675.1) with the *lastal* program in

180    *LAST* (Kiełbasa et al. 2011; Frith & Kawaguchi 2015) using parameters optimized for distantly

181    related species (*-m100 -E0.05*). Given the deep divergence between the two species, we kept

182    low-confidence alignments (*last-split -m1*). We filtered alignments shorter than 500 bp and

183    visualized syntenic relationships only for silverside scaffolds longer than 1 Mb ('chromosome

184    assembly', see below) using the software CIRCA (omgenomics.com/circa).

185

186    *Repeat and gene annotation*

187    We annotated the Atlantic silverside genome using a combination of the *BRAKER2* (Hoff et al.

188    2019) and *MAKER* (Holt & Yandell 2011) pipelines, which combine repeat masking, *ab initio*

189    gene predictor models and protein and transcript evidence for *de novo* identification and

190    annotation of genes. To annotate repetitive elements, we first identified repeats *de novo* in the

191    Atlantic silverside genome using *Repeatmodeler* (Smit & Hubley 2008) and NCBI as a search

192    engine and combined the resulting species-specific library with a library of known repeats in

193    teleosts (downloaded from the RepBase website (Bao et al. 2015) in July 2018). The merged

194    libraries were then used to annotate repeats in the Atlantic silverside genome with *Repeatmasker*

195    (Smit et al. 2015). We then filtered annotated repeats to only keep complex repeats for soft-

196    masking. Next, we used *BRAKER2* to train *AUGUSTUS* (Stanke et al. 2006; Stanke et al. 2008;

197    Buchfink et al. 2015) on the soft-masked genome with unpublished mRNA-seq evidence from 24

198    Atlantic silverside individuals from different populations and developmental stages, along with

9

199    protein homology evidence from six different teleost species (medaka [*Oryzias latipes*], tilapia

200    [*Oreochromis aureus*], platyfish [*Xiphophorus maculatus*], zebrafish [*Danio rerio*], stickleback

201    [*Gasterosteus aculeatus*] and fugu [*Takifugu rubripes*]), which were downloaded from

202    ensemble.org (Ensembl 98; Cunningham et al. 2019) and the UniProtKB (Swiss-Prot) protein

203    database. Second, we ran five rounds of annotation in *MAKER* using different input datasets. The

204    first round of *MAKER* was performed on the genome with only complex repeats masked using

205    the non-redundant transcriptome of Atlantic silverside (Therkildsen and Palumbi 2017,

206    Therkildsen and Baumann 2020) as mRNA-seq evidence, and the six protein sequence datasets

207    from other species as protein homology evidence. We then trained *SNAP* (Korf 2004) on the

208    output of the initial *MAKER* run for *ab initio* gene model prediction. We ran *MAKER* a second

209    time adding the SNAP *ab initio* gene predictions. Using the *MAKER* output from this second

210    round, we re-trained *SNAP* and ran *MAKER* three additional times (round 3 to 5) including the

211    updated *SNAP* gene predictions, the *AUGUSTUS* gene predictions from *BRAKER2* and the

212    updated *MAKER* annotation.

213        Lastly, we performed a functional annotation using *Blast2GO* in *Omnibox v.1.2.4* (Götz

214    et al. 2008) utilizing the UniProtKB (Swiss-Prot) database and *InterProScan2* results. Annotated

215    Atlantic silverside nucleotide sequences for all predicted genes were blasted against the

216    UniProtKB database using *DIAMOND* v. 0.9.34 (Buchfink et al. 2015) with an e-value cutoff of

217    $10^{-5}$. *InterProScan2* was used to annotate proteins with *PFAM* and *Panther* annotations and

218    identify GO terms. *Blast2GO* default mapping and annotation steps were performed using both

219    lines of evidence to create an integrated annotation file.

220

221

10

222    *Comparison of sequence and structural standing genetic variation between populations*

223    As Atlantic silversides from Georgia show strong genomic differentiation from populations

224    further north, primarily concentrated in large haploblocks on four chromosomes (Therkildsen et

225    al. 2019; Wilder et al. 2020), we also sequenced the genome of a representative individual from

226    Mumford Cove, Connecticut, USA (N 41.32°, W 72.02°) collected in June 2016 for comparison.

227    Genomic DNA was extracted from muscle tissue using the DNeasy Blood and Tissue kit

228    (Qiagen) and normalized to 40 ng/µl. We prepared a genomic DNA library using the TruSeq

229    DNA PCR-free library kit (Illumina) following the manufacturer's protocol for 550 bp insert

230    libraries. The shotgun library was sequenced using paired-end 150 bp reads on an Illumina

231    HiSeq4000.

232        We estimated genome size and heterozygosity from the raw data from this shotgun

233    library using the same k-mer approach as for the Georgia individual described above. To

234    compare our heterozygosity estimates in Atlantic silversides from Connecticut and Georgia with

235    other fish species, we searched the literature for heterozygosity estimates from Genomescope

236    with the keywords "Genomescope heterozygosity fish", or from variant calling methods in other

237    fish genomes, using Google Scholar. We also estimated heterozygosity directly by calculating

238    the proportion of heterozygous sites in each genome. For the Georgia individual we used the

239    processed 10X data as above. For the Connecticut individual we trimmed adapters and low-

240    quality data from the raw shotgun data in *Trimmomatic* (Bolger et al. 2014). We mapped data

241    from the two libraries to the chromosome assembly (only the largest 27 scaffolds - see Results)

242    with *bwa mem* (Li & Durbin 2009) and removed duplicates with *samblaster* (Faust & Hall 2014).

243    We called variants with *bcftools mpileup* and *bcftools call* (Danecek et al. 2014). As areas of the

244    genome covered by more than twice the mean sequencing depth could represent repetitive areas

11

245    or assembly artefact, we calculated genome coverage for each of the two libraries with

246    *genomeCoverageBed* from *BEDtools* (Quinlan & Hall 2010) and identified the depth mode from

247    the calculated distribution (95x for the southern genome and 74x for the northern genome). We

248    then filtered variants that were flagged as low-quality, that had mapping quality below 20,

249    sequencing depth below 20, and more than twice the mode sequencing depth for each of the two

250    libraries using *bcftools filter* (Li et al. 2009). To accurately estimate the proportion of

251    heterozygous sites in the genome, we subtracted the number of sites that had sequencing depth

252    below 20 and above twice the mode sequencing depth from the total genome size (to get the sum

253    of sites that could be identified as either homozygous or heterozygous based on our criteria). To

254    visualize variation along the genome, we plotted estimates of heterozygosity in 50-kb sliding

255    windows along the genome for each of the two individuals using the *qqman* package (Turner

256    2014) in R (R Core Team 2019). To assess the reduction in diversity in protein-coding regions

257    due to positive and purifying selection, we calculated heterozygosity in the regions annotated as

258    coding sequences only and compared this to the genome-wide estimate.

259          Finally, we identified structural variants (SVs) segregating between the Connecticut and

260    Georgia genomes using *Delly2 v.0.8.1* (Rausch et al. 2012). For this analysis we used the

261    shotgun library data (74x coverage) from Connecticut mapped to the Georgia reference genome

262    as described above. We called SVs using the command *delly call* and default settings. As

263    genotyping a single individual in *Delly* is prone to false positives we applied the following

264    stringent filters: We retained only homozygous SVs (*vac=2*) that passed quality filters (*PASS*)

265    and that had at least 20 reads supporting the variant calls, whether they came from paired-end

266    clustering or split-read analysis or a combination of the two, but not more than 100 reads since

267    these could be due to repetitive elements in the genome. As *Delly2* outputted redundant

12

268    genotypes, e.g. inversions that had slightly different breakpoints were reported as independent

269    variants, we used *bedtools merge* to merge these overlapping features. To validate duplication

270    calls we also calculated coverage for each of these variants and retained only those putative

271    duplications that had coverage more than 1.8-fold the whole genome sequencing depth (74x).

272         To confirm the large SVs observed between the two genomes examined, we generated a

273    second Hi-C library from an Atlantic silverside individual caught in Mumford Cove, Connecticut

274    in June 2016 (different from the sample used for the shotgun assembly). Liver tissue was excised

275    and digested for 2 hours in collagenase digestion buffer (perfusion buffer plus 12.5 μM CaCl2

276    plus collagenases II and IV (5 mg/ml each)). The cell suspension was then strained through a 100

277    μm cell strainer, washed with 1 ml cold PBS three times, resuspended in 45 ml PBS, and

278    quantified in a hemocytometer. The cross-linking protocol was modified from Belton et al.

279    (2012) as follows. 1.25 ml of 37% formaldehyde was added twice to the cell preparation, then

280    incubated at room temperature for 10 minutes, inverting every 1-2 minutes. To quench the

281    formaldehyde in the reaction, 2.5 ml of 2.5 M glycine was added three times. The sample was

282    incubated at room temperature for 5 minutes, then on ice for 15 minutes to stop the cross-linking.

283    The cells were pelleted by centrifugation (800g for 10 min), and the supernatant was removed.

284    The sample thus obtained was flash frozen in liquid nitrogen and stored at -80°C. Hi-C library

285    preparation was performed as described previously (Belaghzal et al. 2017), except that ligated

286    DNA size selection was omitted. 50 million fish liver cells were digested with *DpnII* at 37°C

287    overnight. DNA ends were filled with biotin-14-dATP at 23°C for 4 hours. DNA was then

288    ligated with T4 DNA ligase at 16°C overnight. Proteins were removed by treating ligated DNA

289    with proteinase-K at 65°C overnight. Purified, proximally ligated molecules were sonicated to

290    obtain an average fragment size of 200 bp. After DNA end repair, dA-tailing and biotin pull

291    down, DNA molecules were ligated to Illumina TruSeq sequencing adapters at room temperature

292    for 2 hours. Finally, the library was PCR-amplified and finalized following the Illumina TruSeq

293    Nano DNA Sample Prep kit manual. Paired-end 50 bp sequencing was performed on a

294    HiSeq4000.

295          The two Hi-C libraries from Connecticut and Georgia (the latter prepared by Dovetail)

296    were mapped to the Atlantic silverside chromosome assembly using the *Distiller* pipeline

297    (github.com/mirnylab/distiller-nf). Interaction matrices were binned at 50 and 100 kb resolution

298    and intrinsic biases were removed using the Iterative Correction and Eigenvector decomposition

299    (ICE) method (Imakaev et al. 2012). Large inversions (> 1 Mb) were identified by visual

300    inspection of Hi-C maps as discontinuities that would be resolved when the corresponding

301    section of the chromosomes were to be inverted (Dixon et al. 2018; Corbett-Detig et al. 2019).

302    These discontinuities generate a distinct "butterfly pattern" with signals of more frequent Hi-C

303    interactions where the projected coordinates of the breakpoints meet.

304

**Results**

*Genome assembly and assessment of completeness*

307    We obtained the best draft assembly (with the highest contiguity; N50 = 1.3 Mb) from the 10X

308    data when we used 270 million reads as input to *Supernova*. Contiguity increased more than 2-

309    fold with Dovetail Chicago data (scaffold N50 = 2.9 Mb) and more than 10-fold with Dovetail

310    Hi-C data (scaffold N50 = 18.2 Mb). Summary statistics for each of the intermediate genome

311    assemblies (10X, Dovetail Chicago, and Dovetail Hi-C) are presented in Table 1. The final

312    assembly – including scaffolds longer than 1 kb only – was 620 Mb in total length. Overall, this

313    assembly showed high contiguity, high completeness and a low proportion of gaps (Table 1).

14

314    Analysis of the presence of BUSCO genes showed that only 5.9% of the Actinopterygii gene set

315    were missing from the assembly. Although the number of missing genes did not decrease

316    dramatically from the 10X assembly to the final assembly (from 6.6 to 5.9%), the addition of

317    proximity ligation data (Chicago and Hi-C) increased the number of complete genes (from 88.1

318    to 89.6%) and decreased the number of duplicated (from 4.1 to 2.9%) and fragmented genes

319    (from 5.3 to 4.5%). Contiguity did not come at the cost of increased gappiness, as stretches of

320    N's made up only 3% of the final assembly. The reduction of the assembly to its longest 27

321    scaffolds ('chromosome assembly'- a 25% reduction in sequence) increased missing genes by

322    only 3.1% and reduced duplicated genes to 1.9%. K-mer analyses based on raw data from the

323    reference genome estimated a genome size of 554 Mb, 76 Mb shorter than the final assembly and

324    88 Mb longer than the chromosome assembly.

325

326    *Synteny with Medaka*

327    The alignment of the 27 largest Atlantic silverside scaffolds to the medaka genome revealed a

328    high degree of synteny conservation, especially considering the evolutionary distance between

329    the two species. Each Atlantic silverside scaffold mapped mostly to only one medaka

330    chromosome, and 22 of the 24 medaka chromosomes had matches with only one Atlantic

331    silverside scaffold each (Fig. 1). Two medaka chromosomes, 1 and 24, had matches with three

332    and two silverside scaffolds, respectively (Fig. 1). Based on these results, karyotype data

333    confirming that the medaka and silverside have the same number of chromosomes (Uwa &

334    Ojima 1981; Warkentine et al. 1987), and additional support from the Hi-C data from the

335    Connecticut individual, we ordered and renamed the Atlantic silverside scaffolds according to

336    the orthologous medaka chromosomes. We grouped the three and two scaffolds that mapped to

15

337   medaka chromosomes 1 and 24, respectively, into one pseudo-chromosome each and renamed

338   them accordingly. Although we did not observe large interchromosomal rearrangements in the

339   alignment of the silverside and medaka genomes (Fig. 1), intrachromosomal rearrangements

340   were common (Fig. 1; Fig. S1). The most conspicuous chromosomal rearrangements were large

341   inversions, intrachromosomal translocations and duplications (Fig. 1; Fig. S1). On chromosomes

342   8, 11, 18 and 24, where large geographically differentiated haploblocks were identified among

343   natural silverside populations, several translocations and inversions were evident, indicating poor

344   intrachromosomal synteny (Fig. 1). This was also the case for most of the other chromosomes

345   (Fig. S1).

346

347   *Repeat and gene annotation*

348   The identified repetitive elements made up 17.73% of the Atlantic silverside genome, in line

349   with expectations based on fish species with similar genome sizes (Yuan et al. 2018). The

350   biggest proportion of these repeats was made up of interspersed repeats (15.34% of the genome),

351   while transposable elements constituted 8.83% of the genome overall (0.90% of SINEs, 2.79%

352   of LINEs, 1.54% of LTR elements, and 3.60% of DNA elements). Our gene prediction pipeline

353   identified a total of 21,644 protein coding genes, a number consistent with annotated gene counts

354   in other fish species (Lehmann et al. 2019; Ozerov et al. 2018). Analysis in *Blast2GO* based on

355   homology and *InterProScan2* resulted in functional annotation of 17,602 out of the 21,644

356   protein coding genes (81.3%; https://github.com/atigano/Menidia_menidia_genome/annotation/).

357   Further, *InterProScan2* detected annotations (*Panther* or *PFAM*) for an additional 1,511 genes,

358   for which no BLAST results were obtained.

359

16

*Sequence and structural standing variation*

K-mer analyses based on raw data resulted in similar estimates of genome sizes and levels of

heterozygosity in the two samples from Georgia and Connecticut: genome size estimates differed

by 20 Mb (554 Mb and 535 Mb in the Georgia and Connecticut individual, respectively) and

heterozygosity estimates differed by 0.09% (1.76% and 1.67% in Georgia and Connecticut,

respectively). Direct estimates of heterozygosity, i.e. based on the number of called heterozygous

sites in the genome, were slightly lower and differed by 0.14% between individuals (1.32% and

1.46% in Georgia and Connecticut, respectively). Together, these estimates concordantly

indicate that standing sequence variation in this species is very high (Kajitani et al. 2014), with 1

in every ~66 bp being heterozygous within each individual. These heterozygosity estimates are

higher than all comparable estimates reported for other fish species, though of similar magnitude

to the European sardine and two eel species (Table 2). Heterozygosity varied substantially across

the genome. Within each chromosome, heterozygosity was highest toward the edges of each

chromosome, presumably in areas corresponding to telomeres, decreased towards the center in a

U-shape fashion, and showed a deep dip in which the number of heterozygous sites approached

zero, consistent with the location of putative centromeres (Fig. 2b). Additionally, the proportions

of variable sites in coding regions was ~50% of whole genome level estimates (0.68% and

0.70% in Georgia and Connecticut, respectively). Swaths of low heterozygosity were particularly

evident on chromosomes 18 and 24, two of four chromosomes with highly differentiated

haploblocks (Fig. 2a,b).

We identified a total of 4,900 SVs - including insertions, deletions, duplications and

inversions (Supplementary File) - between the reference genome generated from Georgia

samples and the re-sequenced individual from Connecticut. *Delly2* indicated that insertions were

17

383 small (42-83 bp) and affected a negligible proportion of the genome, while deletions were larger

384 and more abundant, covering 15% of the genome sequence. As an insertion in one genome

385 corresponds to a deletion in the other genome depending on which individual is used as

386 reference, the discrepancy between insertions and deletions is an artefact of mapping short-read

387 sequences to a single reference, i.e. inserted sequences found only in Connecticut are not present

388 in the reference and thus are not mapped. These results highlight the difficulties in identifying

389 insertions and estimating their sizes from short reads. Our analysis detected a small number of

390 duplications, covering only 0.1% of the genome. In contrast, we identified 662 inversions

391 ranging from 203 bp to 12.6 Mb in size. In total, inversions affected 109 Mb, or 23%, of the

392 reference genome sequence. Twenty-nine inversions were larger than 1 Mb, and five larger than

393 5 Mb (genomic locations in Fig. 2a and in Supplementary File). *Delly2* identified large

394 inversions (> 1 Mb) on all four chromosomes with previously identified haploblocks  . The

395 largest inversion (~12 Mb) was identified on chromosome 8; chromosome 11 had two 1.2-Mb

396 inversions that were 7 Mb apart; chromosome 18 had a 7.4 Mb inversion and chromosome 24

397 had two inversions, the first one spanning 9.4 Mb and followed by another one at a distance of

398 76 kb, spanning 2.3 Mb (Fig. 2a).

399   The independent Hi-C data from Connecticut (which was not used for genome

400 scaffolding) supported a high degree of accuracy in the overall assembly into chromosomes, as

401 indicated by the strong concentration of data points along the diagonal rather than elsewhere in

402 the contact maps (Fig. 3). The contact maps also readily detected large-scale inversions (> 1 Mb)

403 between the individual from Connecticut and the reference assembly from Georgia in three of the

404 four chromosomes with haploblocks, i.e. 8, 18, and 24 (Fig. 3, Supplementary File). The missed

405 detection of the inversions on chromosome 11 could either be due to their relatively smaller

18

406    sizes, barely exceeding the detection threshold from Hi-C data, or because both inversion

407    orientations segregate where the Connecticut individual used for Hi-C was sampled (Wilder et al.

408    2020). The breakpoints of the 12.6 and 9.4 Mb inversions on chromosomes 8 and 24,

409    respectively, matched very closely those identified by *Delly2*, although the second 2.3 Mb

410    inversion on chromosome 24 was not supported by Hi-C data (Figs. 2a, 3, Supplementary File).

411    On chromosome 18, Hi-C data showed a complex series of nested and/or adjacent inversions

412    spanning ~8.8 Mb in total, in contrast with the single inversion, and ~1.3 Mb shorter, identified

413    by *Delly2* (Figs. 2a, 3, Supplementary File). Additional large inversions were detected from the

414    Hi-C data on chromosomes 4, 7 and 19. Of these, the inversion on chromosome 19 was not

415    identified from the analysis of shotgun data with *Delly2*, while those on chromosome 4 and 7

416    were, although with only one matching breakpoint for the inversion on chromosome 4 (Figs. 2a,

417    3, Supplementary File). Note that the identification of SVs from shotgun and Hi-C data were

418    carried out by two different authors, and blindly from each other.

419

420    **Discussion**

421    We generated a chromosome-level assembly of the Atlantic silverside genome by integrating

422    long-range information from synthetic long reads from 10X Genomics, *in vitro* proximity

423    ligation data from Chicago libraries, and Hi-C proximity ligation data from whole cells. The

424    resulting assembly had high contiguity and completeness. Based on karyotype information (Uwa

425    & Ojima 1981; Warkentine et al. 1987), chromosome-level synteny with medaka, and Hi-C maps

426    we reduced the 27 largest scaffolds to 24 putative chromosomes. This chromosome assembly is

427    88 Mb shorter than the genome size estimated through k-mer analysis, but has a lower number of

428    duplicated genes, and only slightly fewer missing genes than the full assembly despite a

19

429    substantial reduction in total sequence. If the proportion of complete genes in the chromosome

430    assembly is, in fact, a good proxy for genome completeness, then the scaffolds that are not

431    placed in chromosomes are mostly sequences that are repetitive, redundant, or that should fill

432    gaps in the assembled chromosomes.

433         Heterozygosity within a sequenced individual can result in alternative alleles getting

434    assembled into distinct scaffolds, even in genomes much less heterozygous than the Atlantic

435    silverside (Kajitani et al. 2014; Tigano et al. 2018), so we expect some redundancy in our

436    assembly. Considering the abundance of SVs between the two sequenced individuals, structural

437    variation also may have contributed to the high number of smaller scaffolds not included in the

438    chromosome assembly, as heterozygous SVs are notoriously hard to assemble (Huddleston et al.

439    2017). Nonetheless, the Atlantic silverside genome adds to the increasing number of high-quality

440    fish reference genome assemblies, with the sixth highest contig N50 (202.88 kb) and the sixth

441    highest proportion of the genome contained in chromosomes (84%, based on the genome size

442    estimate from the k-mer analysis) compared to 27 other chromosome-level fish genome

443    assemblies (Lehmann et al. 2019).

444         Patterns of synteny between the Atlantic silverside and the relatively distantly related

445    medaka are consistent with comparisons among other teleost genomes up to hundreds of millions

446    of years diverged: rearrangements are rare among chromosomes but common within (Amores et

447    al. 2014; Rondeau et al. 2014; Miller et al. 2019; Pettersson et al. 2019). Consistent with this,

448    anchoring Atlantic silverside transcriptome contigs on to medaka genome enabled the

449    identification of four large haploblocks associated with fishery-induced selection in the lab

450    and/or putative adaptive differences in the wild (Therkildsen et al. 2019; Wilder et al. 2020).

451    However, the high degree of intrachromosomal rearrangements between the two species, and

452   generally among teleosts, prevented an accurate characterization of the extent of these

453   haploblocks and the analysis of structural variation. Differentiation between the northern and

454   southern haplotypes seemed to extend across almost the entire length of three of the four

455   chromosomes with haploblocks when data were oriented to medaka (Therkildsen et al. 2019;

456   Wilder et al. 2020). However, the abundant intrachromosomal rearrangements between medaka

457   and Atlantic silverside chromosomes (Fig. 1; Fig. S1), and the detection of large inversions in

458   each of these four chromosomes (Figs. 2a,3) suggest that differentiation is concentrated in, and

459   possibly maintained by, these inversions, which, albeit large, do not span whole chromosomes.

460          Our analysis of two genomes sequenced at high coverage suggested that levels of

461   standing genetic variation, both sequence and structural, are extremely high in the Atlantic

462   silverside. To our knowledge, our estimates of heterozygosity in a single individual are the

463   highest reported for any fish species to date, including those with large census population sizes

464   (Table 2). For example, heterozygosity, which is equivalent to nucleotide diversity ($\pi$) in one

465   individual, in one single Atlantic silverside genome was higher than, or on par with, $\pi$ estimates

466   based on 43-50 individuals of Atlantic killifish, a species considered to have 'extreme' levels of

467   genomic variation with $\pi$ ranging from 0.011 to 0.016 (Reid et al. 2017, 2016). Compared to

468   other vertebrates, genome heterozygosity in the Atlantic silverside was more than double the

469   highest estimate reported for birds (0.7% in the thick-billed murre *Uria lomvia*; Tigano et al.

470   2018) and higher than the population-based 0.6-0.9% estimates in the rabbit (*Oryctolagus*

471   *cuniculus*), one of the mammals with the highest genetic diversity (Carneiro et al. 2014). Among

472   a collection of genome-wide $\pi$ estimates - mostly population-based - across 103 animal, plant

473   and fungal populations or species, only three insects and one sponge had $\pi$ estimates higher than

474   the Atlantic silverside (Robinson et al. 2016 and references therein). This unusually high level of

475 standing sequence diversity is likely due to huge population sizes with estimated $N_e$ exceeding

476 100 million individuals (Lou et al. 2018), and may underpin the remarkable degree of adaptive

477 divergence and rapid responses to selection documented for the species.

478   Variation in $\pi$ across the genome has been associated with variation in recombination

479 rates, with higher diversity and recombination rates in smaller chromosomes and in proximity of

480 telomeres in fish, mammals and birds (Ellegren 2010; Murray et al. 2017; Sardell et al. 2018;

481 Tigano et al. 2020). In the Atlantic silverside, the decrease of heterozygosity from the ends

482 towards the center of each chromosome is consistent with decreasing recombination rates as

483 distance from the telomeres increases (Haenel et al. 2018; Sardell et al. 2018). However, in

484 addition to this U-shape pattern, heterozygosity shows a dramatic, narrow dip in each

485 chromosome far from the center of chromosomes, suggesting a strong centromere effect.

486 Although striking differences exist between sexes and across taxa, recombination is generally

487 reduced or suppressed around centromeres (Sardell & Kirkpatrick 2020). The Atlantic silverside

488 karyotype, with only four metacentric and 20 non-metacentric chromosomes (i.e. submetacentric,

489 subacrocentric, and acrocentric; Warkentine et al. 1987), further supports that these dips in

490 heterozygosity are associated with centromeres, as the non-metacentric chromosomes enable the

491 distinction between the effect of centromeres from the effect of distance from telomeres. In

492 forthcoming work, linkage mapping will allow us to quantify the relative effects of centromeres

493 and telomeres on local recombination rates and ascertain whether the recombination landscape is

494 different between sexes.

495   We report a 50% reduction in heterozygosity in coding sequences compared to whole

496 genome estimates, confirming the expectation that estimates based on exome data are not

497 representative of whole-genome levels of standing variation. Even though the magnitude of the

22

498    reduction in $\pi$ within coding regions is similar to levels reported in the Atlantic killifish (Reid et

499    al. 2017) and in the butterfly *Heliconius melpomene* (Martin et al. 2016), a substantially greater

500    reduction is seen in the collared flycatcher (86%; Dutoit et al. 2017), suggesting that the

501    distribution of diversity in a genome, including the difference between coding and non-coding

502    sequence, is likely idiosyncratic to the population or species examined. Once again, a paucity of

503    data from other species prevents us from making generalizations or identifying differences on the

504    expected reduction in diversity in coding compared to non-coding regions across taxa, while at

505    the same time it highlights the importance of estimating and reporting basic diversity statistics

506    for whole genome assemblies.

507        We identified 4,900 structural variants that survived the stringent filters applied to

508    maximize confidence in the identified SVs and to minimize the number of false positives due to

509    genotyping one individual only. Our estimates are likely conservative when we consider that we

510    filtered out all heterozygous SVs, that many SVs, particularly complex ones, are hard to identify

511    or characterize (Chaisson et al. 2019), and that we analyzed only two genomes. Nonetheless, our

512    analyses based on shotgun data show that SVs are abundant, affect a large proportion of the

513    genome, with inversions covering up to 23% of the genome sequence, and range in size from

514    small (< 50 bp) to longer than 10 Mb, with many of the largest inversions further supported by

515    independent Hi-C data. Sunflower species of the genus *Helianthus* show a similar proportion of

516    sequence covered by inversions (22%; Barb et al. 2014), although these were detected in

517    comparisons between species (1.5 million years diverged) rather than within species. The few

518    studies available on other species show that structural variation tends to affect a larger portion of

519    the genome than single nucleotide polymorphisms (SNPs), but in proportions far lower than what

520    we report here for the Atlantic silverside. For example, structural variation, including indels,

23

521    duplication and inversions, covered three times more bases than SNPs did across six individuals

522    of Australian snapper (*Chrysophrys auratus*; Catanach et al. 2019); short indels alone affected

523    4% of the genome of two individuals from the same population in the cactus mouse (*Peromyscus*

524    *eremicus*; Tigano et al. 2020); inversions, duplications and deletions combined affected 3.6% of

525    the genome across 20 individuals of *Tinema* stick insects (Lucek et al. 2019); and in cod (*Gadus*

526    *morhua*) inversions covered ~7.7% of the genome (Wellenreuther & Bernatchez 2018 and

527    references therein). Although levels of structural variation in the Atlantic silverside are extreme

528    in comparison to these studies, a direct comparison with these and other species is hampered by a

529    paucity of data and lack of common best practices for SVs genotyping (Mérot et al. 2020):

530    differences in sampling, approaches, data types and filtering prevent comparisons similar to

531    those made for standing sequence variation here and in other studies (Corbett-Detig et al. 2015;

532    Robinson et al. 2016). Given the fast rate at which high-quality reference genomes are now

533    generated, this will hopefully start to change.

534         The simple and affordable strategy we adopted only requires sequencing of a single

535    additional shotgun library prepared from a second individual - possibly from a differentiated

536    population to capture a broader representation of intraspecific variation - and could be easily

537    applied in other studies to start describing variation in the prevalence and genome coverage of

538    SVs across taxa. Here, an additional Hi-C library then allowed us to discover that the putative

539    inversion on chromosome 18 was larger than indicated by the analysis of shotgun data and was

540    actually constituted by a combination of two or more nested inversions. The apparent

541    discrepancy between the breakpoints of the largest inversions identified using the two data types

542    could reflect biological variation between the individuals analyzed. Alternatively, they may be

543    caused by the different strengths and limitations of the underlying analytical approaches,

544    including the fact that the identification of SVs was computational from shotgun data, while it

545    was manually curated from Hi-C data. Although the analysis of only two individuals does not

546    capture the full spectrum of intra- and inter-population variation, integrating different approaches

547    has allowed us to identify a set of high-confidence SVs to be validated and genotyped in a larger

548    number of individuals with lower coverage data (Mérot et al. 2020).

549        The joint analysis of sequence and structural variation reveals interesting features of the

550    previously identified haploblocks. The chromosome-level assembly of the Atlantic silverside

551    genome a) confirms that previously identified large haploblocks (Wilder et al. 2020) are

552    associated with inversions and allows to measure their real extent ; and b) highlights how

553    genomic heterogeneity is multidimensional by revealing that even haploblocks showing similar

554    patterns of differentiation can show vastly different patterns of genetic diversity. On

555    chromosomes 18 and 24, large swaths of reduced heterozygosity (Fig. 2b) are associated with an

556    inversion affecting the same area, which strongly indicates that the inversion promotes

557    differentiation between genomes from Connecticut and Georgia in this region, likely through

558    suppressed recombination. Of note, however, the segment of chromosome 24 preceding the

559    inversion (0-722 kb) shows an even stronger reduction in heterozygosity than the adjacent

560    inversion. While this additional reduction may be due to stronger recombination suppression in

561    this area, the mechanism explaining this pattern remains to be investigated. In contrast, no

562    reduction in diversity is associated with the inversion on chromosome 8 - the largest of them all

563    (12.6 Mb) - or with the smaller inversions on chromosome 11.  Such differences among

564    haploblocks likely reflect idiosyncratic evolutionary histories and adaptive significance of the

565    underlying inversions, whose investigation is now enabled by the chromosome-level genome

566    assembly that we presented here. Hence, our analyses provide an empirical example of the

25

567    importance of analyzing both sequence and structural variation to understand the mechanism

568    underpinning the heterogeneous landscape of genomic diversity and differentiation.

569         Building on prior analysis based on in silico exome capture (Therkildsen & Palumbi

570    2017; Therkildsen et al. 2019; Therkildsen & Baumann 2020), this newly assembled reference

571    genome provides an important resource for using the Atlantic silverside as a powerful model for

572    investigating many outstanding questions in adaptation genomics, for example related to the

573    abundance, distribution and adaptive value of structural variants; the relative role of coding and

574    non-coding regions; the importance of sequence variation vs. structural variation in both human-

575    induced evolution and local adaptation; and the demographic and evolutionary factors generating

576    the genomic landscape of diversity and differentiation in this and other species.

577

578    **Acknowledgements**

587

588

589

26

590    **Author contributions**

591    AT and NOT designed the study with input from JD; AJ performed the gene annotation; AW

592    collected samples and performed lab work; AT, NOT, AW, YZ, AN and JD generated and

593    analyzed the data; NOT and JD funded the project. AT wrote the paper with critical input from

594    all authors.

595

596    **Data accessibility**

597    The genome assembly and associated sequence data from Georgia and the raw data from the

598    shotgun library from Connecticut will be available under ENA accession number #######.

599    Scripts for the genome assembly and all other analyses can be found at

600    https://github.com/atigano/Menidia_menidia_genome/.

601

602   **Tables and Figures**

603   Table 1. Summary statistics for each of the intermediate and final assemblies produced.

|  | 10X | Dovetail Chicago | Dovetail Hi-C | Final assembly | Chromosome assembly* |
|---|---|---|---|---|---|
| Total length | 645.45 Mb | 647.32 Mb | 647.39 Mb | 620.04 Mb | 465.69 Mb |
| Longest Scaffold | 12,248,921 bp | 12,871,938 bp | 26,678,928 bp | 26,678,928 bp | 26,678,928 bp |
| Number of scaffolds | 99,541 | 80,990 | 80,312 | 42,220 | 27 |
| Number of scaffolds > 1kb | 61,451 | 42,898 | 42,220 | 42,220 | 27 |
| Contig N50 | 39.55 kb | 39.51 kb | 39.51 kb | 105.76 kb | 202.88 kb |
| Scaffold L50/N50 | 83/1.328 Mb | 42/2.936 Mb | 16/18.159 Mb | 15/18.199 Mb | 11/19.68 Mb |
| % gaps | 2.69% | 2.97% | 2.98% | 3.08% | 3.00% |
| BUSCOs** (n=4584) | C:88.1%, F:5.3%, M:6.6% | C:89.5%, F:4.6%, M:5.9% | C:89.6%, F:4.8%, M:5.6% | C:89.6%, F:4.5%, M:5.9% | C:88.3%, F:2.7%, M:9.0% |

604   * The 'chromosome assembly' is the subset of scaffolds > 1 Mb from the 'Final assembly'

605   ** [C=complete, F=fragmented, M=missing]

606

607

28

608 Table 2. Examples of heterozygosity levels in single fish genomes, estimated either with

609 GenomeScope from raw sequencing data or through direct calling of heterozygous sites.

| Common name | Scientific name | Heterozygosity | Method | Reference |
|---|---|---|---|---|
| **Atlantic silverside** | ***Menidia menidia*** | **1.67-1.76%** | GenomeScope | **This study** |
| European sardine | *Sardina pilchardus* | 1.60–1.75% | GenomeScope | Machado et al. 2018 |
| American eel | *Anguilla rostrata* | 1.5-1.6% | GenomeScope | Jansen et al. 2017 |
| European eel | *Anguilla anguilla* | 1.48-1.59% | GenomeScope | Jansen et al. 2017 |
| Pearlscale pygmy angelfish | *Centropyge vrolikii* | 1.36% | GenomeScope | Fernandez-Silva et al. 2018 |
| Marine medaka | *Oryzias melastigma* | 1.19% | GenomeScope | Kim et al. 2018 |
| Large yellow croaker | *Larimichthys crocea* | 1.06% | GenomeScope | Mu et al. 2018 |
| Javafish medaka | *Oryzias javanicus* | 0.96% | GenomeScope | Takehana et al. 2020 |
| Greater amberjack | *Seriola dumerili* | 0.65% | GenomeScope | Sarropoulou et al. 2017 |
| Clownfish | *Amphiprion ocellaris* | 0.60% | GenomeScope | Tan et al. 2018 |
| Hilsa shad | *Tenualosa ilisha* | 0.58-0.66% | GenomeScope | Mollah et al. 2019 |
| Whitefish | *Coregonus sp. "Balchen"* | 0.44% | GenomeScope | De-Kayne et al. 2020 |
| Corkwing wrasse | *Symphodus melops* | 0.40% | GenomeScope | Mattingsdal et al. 2018 |
| Herring | *Clupea harengus* | 0.32% | Variant calling | Martinez Barrio et al. 2016 |
| Golden pompano | *Trachinotus ovatus* | 0.31% | GenomeScope | Zhang et al. 2019 |
| Coelacanth | *Latimeria chalumnae* | 0.28% | Variant calling | Amemiya et al. 2013 |
| NA | *Lucifuga gibarensis* | 0.26% | GenomeScope | Policarpo et al. 2020 |
| Eurasian perch | *Perca fluviatilis* | 0.24–0.28% | GenomeScope | Ozerov et al. 2018 |
| Atlantic cod | *Gadus morhua* | 0.20% | Variant calling | Star et al. 2011 |
| Big-eye mandarin Fish | *Siniperca knerii* | 0.16% | GenomeScope | Lu et al. 2020 |
| Threespine stickleback | *Gasteosteus aculeatus* | 0.14% | Variant calling | Jones et al. 2012 |
| Pikeperch | *Sander lucioperca* | 0.14% | GenomeScope | Nguinkal et al. 2019 |
| African arowana | *Heterotis niloticus* | 0.13% | GenomeScope | Hao et al. 2020 |
| Orange clownfish | *Amphiprion percula* | 0.12% | GenomeScope | Lehmann et al. 2019 |
| Murray cod | *Maccullochella peelii* | 0.10% | GenomeScope | Austin et al. 2017 |
| Toothed Cuban cusk-eel | *Lucifuga dentata* | 0.10% | GenomeScope | Policarpo et al. 2020 |

610

611

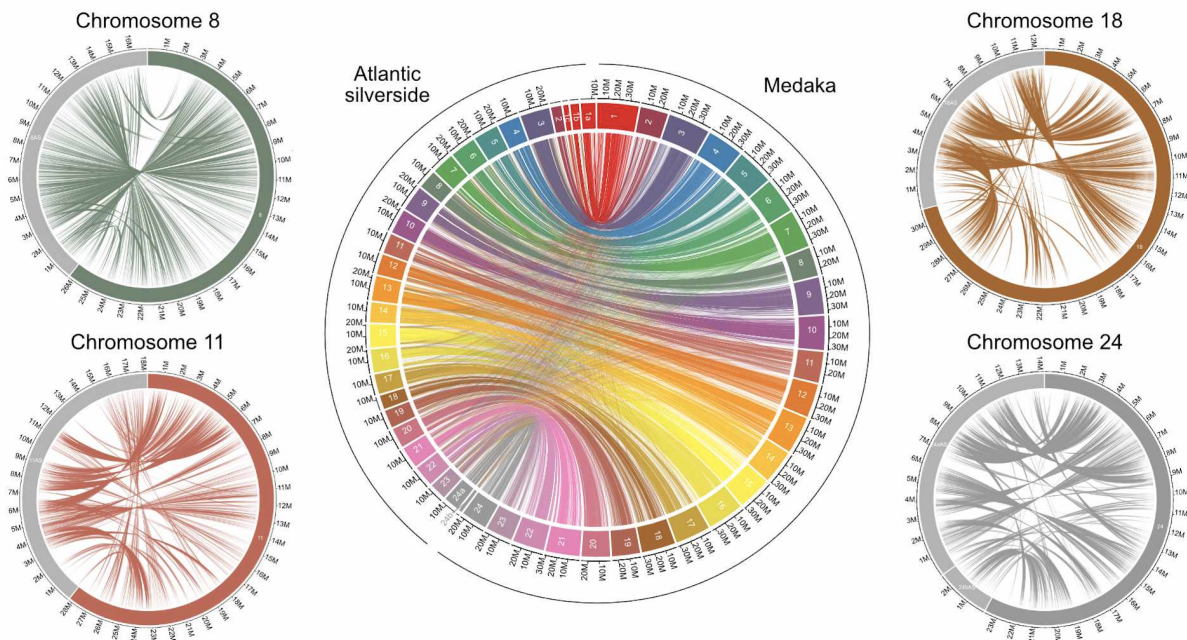612     Table 3. Summary of intraspecific structural variants identified in the Atlantic silverside, and

613     their features.

| SV type | Number of variants | Size range (bp) | Sequence affected (kb) | % genome affected |
|---|---|---|---|---|
| Insertions | 299 | 42-83 | 18 | <0.01% |
| Deletions | 3905 | 38-9,740,501 | 71,754 | 15% |
| Duplications | 34 | 110-150,263 | 479 | 0.1% |
| Inversions | 662 | 203-12,585,625 | 109,201 | 23% |

614

615

616    Figure 1. Circos plots showing synteny between the Atlantic silverside and medaka across all

617    chromosomes in the middle and in the four chromosomes with large haploblocks on the sides.

618    Chromosomes are color-coded consistently among plots and the colored portion of the smaller

619    plots refer to the medaka sequences, while the grey portion to the Atlantic silverside sequences.

620    Alignments shorter than 500 bp were excluded. Fig. S1 shows plots for the remaining

621    chromosomes. Note that the consistently shorter length of the Atlantic silverside genome is

622    consistent with a lower overall estimate of genome size (554 Mb based on k-mer analysis

623    compared to the 700 Mb of the assembled medaka genome). The three and two scaffolds making

624    up chromosomes 1 and 24, respectively, are represented separately here and denoted by small
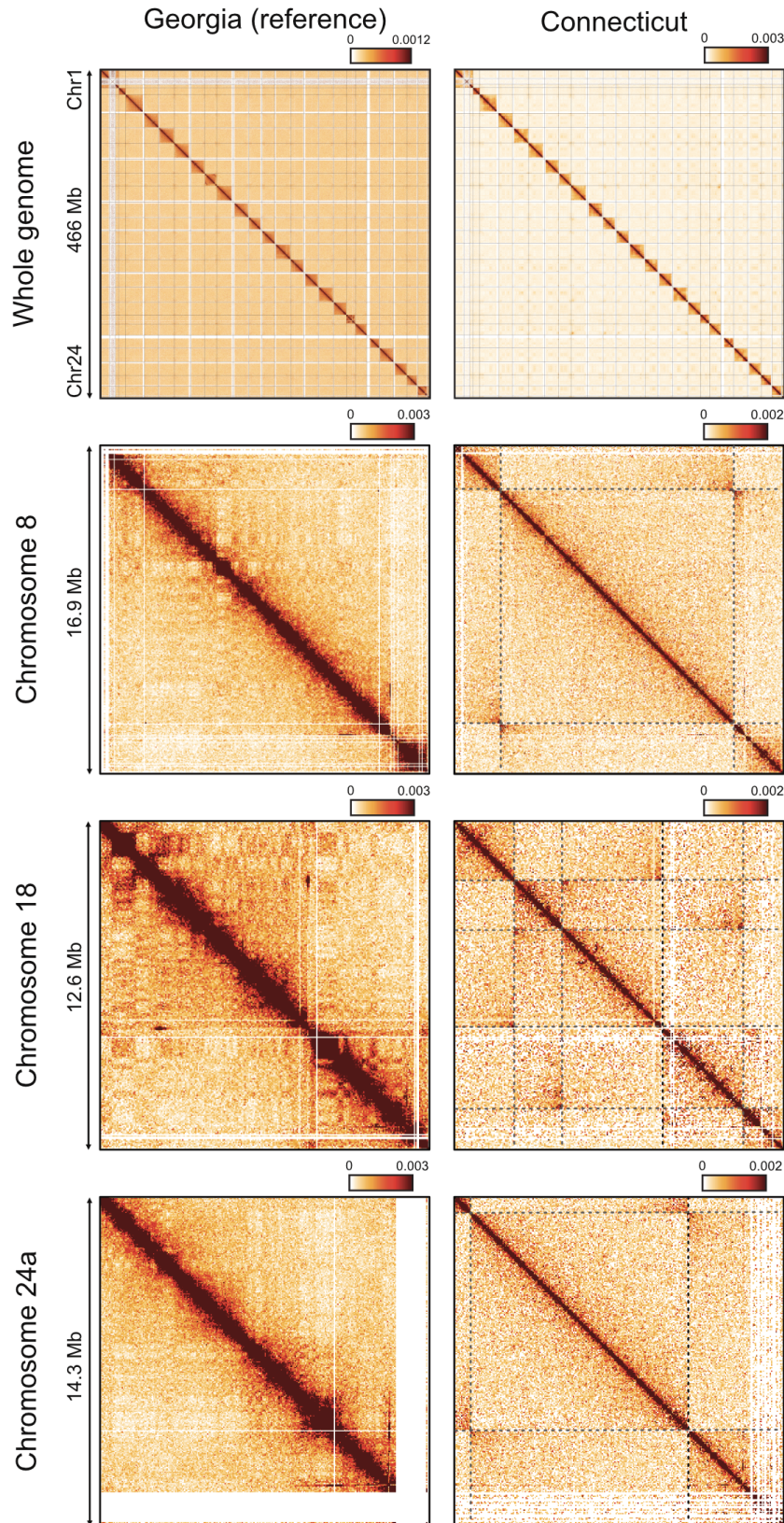
625    letters.



626

627

628

629

31

630    Figure 2. The genomic landscape of structural and sequence variation in Connecticut and

631    Georgia. a) Panel showing large inversions (> 1 Mb) as identified from shotgun and Hi-C data

632    from an individual from Connecticut mapped to the reference genome from Georgia. b)

633    Manhattan plots showing the genomic landscape of variation in heterozygosity in 50 kb moving

634    windows across single genomes from Connecticut and Georgia. The three and two scaffolds

635    making up chromosomes 1 and 24, respectively, are represented separately here and denoted by

636    small letters (e.g., 1a and 24a).



637

638

639

640

641

642

643

644 Figure 3. Hi-C contact maps of data mapped to the chromosome assembly from Georgia. Maps

645 on the left show Hi-C data obtained from the same Georgia individual used to generate the

646 reference assembly (mapped to self), maps on the right show data obtained from a Connecticut

647 individual. Maps in the top panel show data for all the chromosomes binned in 100 kb sections.

648 The three lower panels show data binned in 50 kb sections from each of the three chromosomes

649 showing both large haploblocks in Wilder et al. (2020) and evidence for the presence of

650 inversions from Hi-C data. Dark shades on the diagonal are indicative of high structural

651 similarity between the reference and the Hi-C library analyzed. Dashed lines represent putative

652 inversion breakpoints. The "butterfly pattern" of contacts observed at the point when the dashed

653 lines meet is diagnostic of inversions.

654

34

## References

655

656  Amemiya CT et al. 2013. The African coelacanth genome provides insights into tetrapod
657  evolution. Nature. 496:311–316.

658  Amores A et al. 2014. A RAD-tag genetic map for the platyfish (*Xiphophorus maculatus*) reveals
659  mechanisms of karyotype evolution among teleost fish. Genetics. 197:625–641.

660  Austin CM, Tan MH, Harrison KA, Lee YP, Croft LJ, Sunnucks P, Pavlova A, Gan HM. 2017.
661  De novo genome assembly and annotation of Australia's largest freshwater fish, the Murray cod
662  (*Maccullochella peelii*), from Illumina and Nanopore sequencing read. GigaScience. 6.

663  Bao W, Kojima KK, Kohany O. 2015. Repbase Update, a database of repetitive elements in
664  eukaryotic genomes. Mobile DNA. 6:11.

665  Barb JG, Bowers JE, Renaut S, Rey JI, Knapp SK, Rieseberh LH, Burke JM. 2014.
666  Chromosomal evolution and patterns of introgression in helianthus. Genetics. 197:969–979.

667  Barrett RDH, Schluter D. 2008. Adaptation from standing genetic variation. Trends Ecol. Evol.
668  23:38–44.

669  Belaghzal H, Dekker J, Gibcus JH. 2017. Hi-C 2.0: An optimized Hi-C procedure for high-
670  resolution genome-wide mapping of chromosome conformation. Methods. 123:56–65.

671  Belton J-M, McCord RP, Gibcus, Naumova, Zhan Y, Dekker. 2012. Hi–C: A comprehensive
672  technique to capture the conformation of genomes. Methods. 58:268–276.

673  Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence
674  data. Bioinformatics. 30:2114–2120.

675  Buchfink B, Xie C, Huson DH. 2015. Fast and sensitive protein alignment using DIAMOND.
676  Nat. Methods. 12:59–60.

677  Campagna L, Repenning M, Silveira LF, Fontana CS, Tubaro L. Pablo, and IJ Lovette. 2017.
678  Repeated divergent selection on pigmentation genes in a rapid finch radiation. Science
679  Advances. 3:e1602404.

680  Carneiro M, Rubin C-J, Di Palma F, Albert FW, Alföldi J, Martinez Barrio A, Pielberg G, Rafati
681  N, Sayyab S, Turner-Maier J et al. 2014. Rabbit genome analysis reveals a polygenic basis for
682  phenotypic change during domestication. Science. 345:1074–1079.

683  Catanach A, Crowhurst R, Deng C, David C, Bernatchez L, Wellenreuther M. 2019. The
684  genomic pool of standing structural variation outnumbers single nucleotide polymorphism by
685  threefold in the marine teleost *Chrysophrys auratus*. Mol. Ecol. 28:1210–1223.

686  Chaisson MJP, Sanders AS, Zhao X, Malhotra D, Porubsky D, Rausch T, Gardner EJ, Rodriguez
687  OL, Guo L, Collins RL, et al. 2019. Multi-platform discovery of haplotype-resolved structural
688  variation in human genomes. Nat. Commun. 10:1784.

689  Conover DO, Munch SB. 2002. Sustaining Fisheries Yields Over Evolutionary Time Scales.

690    Science. 297:94–96.

691    Conover DO, Arnott SA, Walsh MR, Munch SB. 2005. Darwinian fishery science: lessons from
692    the Atlantic silverside (*Menidia menidia*). Canadian Journal of Fisheries and Aquatic Sciences.
693    62:730–737.

694    Conover DO, Kynard BE. 1981. Environmental sex determination: interaction of temperature
695    and genotype in a fish. Science. 213:577–579.

696    Conover DO, Present TMC. 1990. Countergradient variation in growth rate: compensation for
697    length of the growing season among Atlantic silversides from different latitudes. Oecologia.
698    83:316–324.

699    Corbett-Detig RB, Said, I, Calzetta M, Gdenetti M, McBroome J, Maurer MW, Petrarca V, della
700    Torre A, Besansky. 2019. Fine-Mapping Complex Inversion Breakpoints and Investigating
701    Somatic Pairing in the *Anopheles gambiae* Species Complex Using Proximity-Ligation
702    Sequencing. Genetics. 213:1495-1511.

703    Corbett-Detig RB, Hartl DL, Sackton TB. 2015. Natural selection constrains neutral diversity
704    across a wide range of species. PLoS Biol. 13:e1002112.

705    Cunningham F, Achuthan P, Akanni W, Allen J, Amode MR, Armean IM, Bennett R, Bhai J,
706    Billis K, Boddu S et al. 2019. Ensembl 2019. Nucleic Acids Res. 47:D745–D751.

707    Danecek P, Schiffels S, Durbin R. 2014. Multiallelic calling model in bcftools (-m).

708    De-Kayne R, Zoller S, Feulner PGD. 2020. A de novo chromosome-level genome assembly of
709    *Coregonus* sp. 'Balchen': One representative of the Swiss Alpine whitefish radiation. Mol. Ecol.
710    Resour. 20:1093-1109.

711    Dixon JR, Xu J, Dileep V, Zhan Y, Song F, Le VT, Yardimci GG, Chakraborty A, Bann DV,
712    Wang Y, et al. 2018. Integrative detection and analysis of structural variation in cancer genomes.
713    Nat. Genet. 50:1388–1398.

714    Dutoit L, Burri R, Nater A, Mugal CF, Ellegren H. 2017. Genomic distribution and estimation of
715    nucleotide diversity in natural populations: perspectives from the collared flycatcher (Ficedula
716    albicollis) genome. Mol. Ecol. Resour. 17:586–597.

717    Ellegren H. 2010. Evolutionary stasis: the stable chromosomes of birds. Trends Ecol. Evol.
718    25:283–291.

719    Faria R, Chaube P, Morales HE, Larsson T, Lemmon AR, Lemmon EM, Rafajlović M, Panova
720    M, Ravinet M, Johannesson K et al. 2019. Multiple chromosomal rearrangements in a hybrid
721    zone between Littorina saxatilis ecotypes. Mol. Ecol. 28:1375–1393.

722    Faust GG, Hall IM. 2014. SAMBLASTER: fast duplicate marking and structural variant read
723    extraction. Bioinformatics. 30:2503–2505.

724    Fernandez-Silva I, Henderson JB, Rocha LA, Simison WB. 2018. Whole-genome assembly of

725 the coral reef Pearlscale Pygmy Angelfish (*Centropyge vrolikii*). Sci. Rep. 8:1498.

726 Frith MC, Kawaguchi R. 2015. Split-alignment of genomes finds orthologies more accurately.
727 Genome Biol. 16:106.

728 Götz S, García-Gómez JM, Terol J, Williams TD, Nagaraj SH, Nueda MJ, Robles M, Talón M,
729 DopazoJ, Conesa A. 2008. High-throughput functional annotation and data mining with the
730 Blast2GO suite. Nucleic Acids Res. 36:3420–3435.

731 Haenel Q, Laurentino TG, Roesti M, Berner D. 2018. Meta-analysis of chromosome-scale
732 crossover rate variation in eukaryotes and its significance to evolutionary genomics. Mol. Ecol.
733 27:2477–2497.

734 Hao S, Han K, Meng L, Huang X, Shi C, Zhang M, Wang Y, Liu Q, Zhang Y, Seim I et al. 2020.
735 Three genomes of Osteoglossidae shed light on ancient teleost evolution. bioRxiv.
736 2020.01.19.911958. doi: 10.1101/2020.01.19.911958.

737 Hice LA, Duffy TA, Munch SB, Conover DO. 2012. Spatial scale and divergent patterns of
738 variation in adapted traits in the ocean. Ecol. Lett. 15:568–575.

739 Hoff KJ, Lomsadze A, Borodovsky M, Stanke M. 2019. Whole-Genome Annotation with
740 BRAKER. In: Gene Prediction: Methods and Protocols. Kollmar, M, editor. Springer New York:
741 New York, NY pp. 65–95.

742 Holt C, Yandell M. 2011. MAKER2: an annotation pipeline and genome-database management
743 tool for second-generation genome projects. BMC Bioinformatics. 12:491.

744 Huddleston J, Chaisson MJP, Steinberg KM, Warren W, Hoekzema K, Gordon D, Graves-
745 Lindsay, Munson KM, Kronenberg ZN, Vives L, et al. 2017. Discovery and genotyping of
746 structural variation from long-read haploid genome sequence data. Genome Res. 27:677–685.

747 Imakaev Mm Fudenberg G, McCord RP, Naumova N, Goloborodko A, Lajoie BR, Dekker J,
748 Mirny LA. 2012. Iterative correction of Hi-C data reveals hallmarks of chromosome
749 organization. Nat. Methods. 9:999–1003.

750 Jansen HJ, Liem M, Jong-Raadsen SA, Dufour S, Weltzien F-A, Swinkerls W, Koelewijn A,
751 Palstra AP, Pelster B, Spaink HP et al. 2017. Rapid de novo assembly of the European eel
752 genome from nanopore sequencing reads. Sci. Rep. 7:7213.

753 Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, Swofford R, Pirun M, Zody
754 MC, White S et al. 2012. The genomic basis of adaptive evolution in threespine sticklebacks.
755 Nature. 484:55–61.

756 Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada M,
757 Nagayasu E, Maruyama H et al. 2014. Efficient de novo assembly of highly heterozygous
758 genomes from whole-genome shotgun short reads. Genome Res. 24:1384–1395.

759 Kess T, Bentzen P, Lehnert SJ, Sylvester EVA, Lien S, Kent MP, Sinclair-Waters M, Morris C,
760 Wringe B et al. 2020. Modular chromosome rearrangements reveal parallel and nonparallel

761    adaptation in a marine fish. Ecol. Evol. 10:638–653.

762    Kiełbasa SM, Wan R, Sato K, Horton P, Frith MC. 2011. Adaptive seeds tame genomic
763    sequence comparison. Genome Res. 21:487–493.

764    Kim H-S, Lee B-Y, Han J, Jeong C-B, Hwang D-S, Lee M-C, Kang H-M, Kim D-H, Lee D, Kim
765    J et al. 2018. The genome of the marine medaka *Oryzias melastigma*. Mol. Ecol. Resour.
766    18:656–665.

767    Korf I. 2004. Gene finding in novel genomes. BMC Bioinformatics. 5:59.

768    Lehmann R, Lightfoot DJ, Schunter C, Mitchell CT, Ohyanagi H, Mineta K, Foret S, Berumen
769    ML, Miller DJ, Aranda M et al. 2019. Finding Nemo's Genes: A chromosome-scale reference
770    assembly of the genome of the orange clownfish *Amphiprion percula*. Mol. Ecol. Resour.
771    19:570–585.

772    Lieberman-Aiden E, van Berkim Nl, Williams L, Imakaev M, Ragoczy T, Telling A, Amit I,
773    Lajoie BR, Sabo PJ, Dorschner MO et al. 2009. Comprehensive mapping of long-range
774    interactions reveals folding principles of the human genome. Science. 326:289–293.

775    Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecais G, Durbin R,
776    1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map format
777    and SAMtools. Bioinformatics. 25:2078–2079.

778    Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
779    Bioinformatics. 25:1754–1760.

780    Lou RN, Fletcher NK, Wylder AP, Conover DO, Therkildsen NO, Searle JB. 2018. Full
781    mitochondrial genome sequences reveal new insights about post-glacial expansion and regional
782    phylogeographic structure in the Atlantic silverside (*Menidia menidia*). Mar. Biol. 165:124.

783    Lucek K, Gompert Z, Nosil P. 2019. The role of structural genomic variants in population
784    differentiation and ecotype formation in *Timema cristinae* walking sticks. Molecular Ecology.
785    28:1224–1237.

786    Lu L, Zhao J, Li C. 2020. High-Quality Genome Assembly and Annotation of the Big-Eye
787    Mandarin Fish (*Siniperca knerii*). G3. 10:877–880.

788    Machado AM, Tørresen OK, Kabeya N, Couto A, Petersen B, Felicio M, Campos PF, Fonseca
789    E, Bandarra N, Lopes-Marques M et al. 2018. 'Out of the Can': A Draft Genome Assembly,
790    Liver Transcriptome, and Nutrigenomics of the European Sardine, *Sardina pilchardus*. Genes.
791    9:485.

792    Marçais G, Kingsford C. 2011. A fast, lock-free approach for efficient parallel counting of
793    occurrences of k-mers. Bioinformatics. 27:764–770.

794    Martinez Barrio A, Lamichhaney S, Fan G, Rafati N, Pettersson M, Zhang H, Dainat J, Ekman
795    D, Höppner M, Jern P et al. 2016. The genetic basis for ecological adaptation of the Atlantic
796    herring revealed by genome sequencing. eLife. 5:e12081

38

797  Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.
798  EMBnet.journal. 17:10–12.

799  Martin SH, Möst M, Palmer WJ, Salazar C, McMillan WO, Jiggins FM, Jiggins CD. 2016.
800  Natural Selection and Genetic Diversity in the Butterfly *Heliconius melpomene*. Genetics.
801  203:525–541.

802  Mattingsdal M, Jentoft S, Tørresen OK, Knutsen H, Hansen MM, Robalo JI, Zagrodzka Z,
803  André C, Gonzalez EB. 2018. A continuous genome assembly of the corkwing wrasse
804  (*Symphodus melops*). Genomics. 110:399–403.

805  Mérot C, Oomen RA, Tigano A, Wellenreuther M. 2020. A Roadmap for Understanding the
806  Evolutionary Significance of Structural Genomic Variation. Trends in Ecology & Evolution.
807  35:561-572

808  Miller JT, Reid NM, Nacci DE, Whitehead A. 2019. Developing a High-Quality Linkage Map
809  for the Atlantic Killifish *Fundulus heteroclitus*. G3. 9:2851–2862.

810  Mollah MBR, Khan MGQ, Islam MS, Alam MS. 2019. First draft genome assembly and
811  identification of SNPs from hilsa shad (*Tenualosa ilisha*) of the Bay of Bengal. F1000Res.
812  8:320.

813  Murray GGR, Soares AER, Novak BJ, Schaefer NK, Cahill JA, Baker AJ, Demboski JR, Doll A,
814  Da Fonseca RR, Fulton TL et al. 2017. Natural selection shaped the rise and fall of passenger
815  pigeon genomic diversity. Science. 358:951–954.

816  Mu Y, Huo J, Guan Y, Fan D, Xiao X, Wei J, Li Q, Mu P, Ao J, Chen X. 2018. An improved
817  genome assembly for *Larimichthys crocea* reveals hepcidin gene expansion with diversified
818  regulation and function. Commun Biol. 1:195.

819  Nguinkal JA, Brunner RM, Verleigh M, Rebi A, los Ríos-Pérez L, Schäfer N, Hadlich F,
820  Stüeken M, Wittenburg D, Goldammer T. 2019. The First Highly Contiguous Genome Assembly
821  of Pikeperch (*Sander lucioperca*), an Emerging Aquaculture Species in Europe. Genes. 10.

822  Ozerov MY, Ahmad F, Gross R, Pukk L, Kahar S, Kisand V, Vasemägi. 2018. Highly
823  Continuous Genome Assembly of Eurasian Perch (*Perca fluviatilis*) Using Linked-Read
824  Sequencing. G3. 8:3737–3743.

825  Pääbo S. 2003. The mosaic that is our genome. Nature. 421:409–412.

826  Pettersson ME, Rochus CM, Han F, Chen J. 2019. A chromosome-level assembly of the Atlantic
827  herring genome—detection of a supergene and other signals of selection. Genome Res. 29:1919-
828  1928

829  Policarpo M, Fumey J, Lafargeas P, Naquin D, Thermes C, Naville M, Dechaud C, Volff J-
830  NCabau C, Klopp C et al. 2020. Contrasted gene decay in subterranean vertebrates: insights from
831  cavefishes and fossorial mammals. bioRxiv.
832  https://www.biorxiv.org/content/10.1101/2020.03.05.978213v1.abstract.

833  Putnam NH, O'Connell BL, Stites JC, Rice BJ, Blanchette M, Calef R, Troll CJ, Fields A,
834  Hartley PD, Sugnet CW et al. 2016. Chromosome-scale shotgun assembly using an in vitro
835  method for long-range linkage. Genome Res. 26:342–350.

836  Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
837  features. Bioinformatics. 26:841–842.

838  Rausch T, Zichener T, Schlattl A, Stütz AM, Benes V, Korbel JO. 2012. DELLY: structural
839  variant discovery by integrated paired-end and split-read analysis. Bioinformatics. 28:i333–i339.

840  Reid NM, Proestou DA, Clark BW, Warren WC, Colbourne JK, Shaw JR, Karchner SI, Hanh
841  ME, Nacci D, Oleksiak MF et al. 2016. The genomic landscape of rapid repeated evolutionary
842  adaptation to toxic pollution in wild fish. Science. 354:1305–1308.

843  Reid NM, Jackson CE, Gilbert D, Minx P, Montague MJ, Hampton TH, Helfrich LW, King BL,
844  Nacci DE, Aluru N et al. 2017. The landscape of extreme genomic variation in the highly
845  adaptable Atlantic killifish. Genome Biol. Evol.

846  Robinson JA, Ortgea-Del Vecchyo D, Fan Z, Kim BY, vonHoldt BM, Marsden CD, Lohmueller
847  KE, Wayne RK. 2016. Genomic Flatlining in the Endangered Island Fox. Curr. Biol. 26:1183–
848  1189.

849  Rondeau EB, Minkley DR, Leong JS, Messmer AM, Jantzen JR, von Schalburg KR, Lemon C,
850  Bird NH, Koop BF. 2014. The genome and linkage map of the northern pike (*Esox lucius*):
851  conserved synteny revealed between the salmonid sister group and the Neoteleostei. PLoS One.
852  9:e102089.

853  Sardell JM, Cheng C, Dagilis AJ, Ishikawa A, Kitano J, Peichel CL, Kirkpatrick M. 2018. Sex
854  Differences in Recombination in Sticklebacks. G3. 8:1971–1983.

855  Sardell JM, Kirkpatrick M. 2020. Sex Differences in the Recombination Landscape. Am. Nat.
856  195:361–379.

857  Sarropoulou E, Sundaram AYM, Kaitetzidou E, Kotoulas G, Gilfillan GD, Papandroulakis N,
858  Mylonas CC, Magoulas A. 2017. Full genome survey and dynamics of gene expression in the
859  greater amberjack *Seriola dumerili*. Gigascience. 6:1–13.

860  Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO:
861  assessing genome assembly and annotation completeness with single-copy orthologs.
862  Bioinformatics. 31:3210–3212.

863  Smit AFA, Hubley R. 2008. RepeatModeler Open-1.0. Available from http://www.
864  repeatmasker. org.

865  Smit AFA, Hubley R, Green P. 2015. RepeatMasker Open-4.0. 2013--2015.

866  Stanke M, Diekhans M, Baertsch R, Haussler D. 2008. Using native and syntenically mapped
867  cDNA alignments to improve de novo gene finding. Bioinformatics. 24:637–644.

868   Stanke M, Schöffmann O, Morgenstern B, Waack S. 2006. Gene prediction in eukaryotes with a
869   generalized hidden Markov model that uses hints from external sources. BMC Bioinformatics.
870   7:62.

871   Star B, Nederbragt AJ, Jentoft S, Grimholt U, Malmstrøm M, Gregers TF, Rounge TB, Paulsen
872   J, Solbakken MH, Sharma A et al. 2011. The genome sequence of Atlantic cod reveals a unique
873   immune system. Nature. 477:207–210.

874   Takehana Y, Zahm M, Cabau C, Klopp C, Roques C, Bouchez O, Donnadieu C, Brrachina C,
875   Journot L, Kawaguchi M, et al. 2020. Genome Sequence of the Euryhaline Javafish Medaka,
876   *Oryzias javanicus*: A Small Aquarium Fish Model for Studies on Adaptation to Salinity. G3.
877   10:907–915.

878   Tan MH, Austin CM, Hammer MP, Lee YP, Croft LJ, Gan HM. 2018. Finding Nemo: hybrid
879   assembly with Oxford Nanopore and Illumina reads greatly improves the clownfish (Amphiprion
880   ocellaris) genome assembly. GigaScience. 7: gix137.

881   Therkildsen NO, Wylder AP, Conover DO, Munch SB, Baumann H, Palumbi SR. 2019.
882   Contrasting genomic shifts underlie parallel phenotypic evolution in response to fishing. Science.
883   365:487–490.

884   Therkildsen NO, Baumann H. 2020. A comprehensive non-redundant reference transcriptome
885   for the Atlantic silverside Menidia menidia. Mar. Genomics. 100738.

886   Therkildsen NO, Palumbi SR. 2017. Practical low-coverage genomewide sequencing of
887   hundreds of individually barcoded samples for population and evolutionary genomics in
888   nonmodel species. Mol. Ecol. Resour. 17:194–208.

889   Tigano A, Colella JP, MacManes MD. 2020. Comparative and population genomics approaches
890   reveal the basis of adaptation to deserts in a small rodent. Mol. Ecol.29:1300-1314.

891   Tigano A, Friesen VL. 2016. Genomics of local adaptation with gene flow. Mol. Ecol. 25:2144–
892   2164.

893   Tigano A, Sackton TB, Friesen VL. 2018. Assembly and RNA-free annotation of highly
894   heterozygous genomes: The case of the thick-billed murre (*Uria lomvia*). Mol. Ecol. Res. 18:79-
895   90

896   Turner SD. 2014. qqman: an R package for visualizing GWAS results using Q-Q and manhattan
897   plots. bioRxiv. 005165. doi: 10.1101/005165.

898   Uwa H, Ojima Y. 1981. Detailed and Banding Karyotype Analyses of the Medaka, *Oryzias
899   latipes* in Cultured Cells. Proc. Jpn. Acad. Ser. B Phys. Biol. Sci. 57:39–43.

900   Van't Hof AE, Campagne P, Rigden DJ, Yung CJ, Lingley J, Quail MA, Hall N, Darby AC,
901   Saccheri IJ. 2016. The industrial melanism mutation in British peppered moths is a transposable
902   element. Nature. 534:102–105.

903   Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowwski, Schatz MC.

904    2017. GenomeScope: fast reference-free genome profiling from short reads. Bioinformatics.
905    33:2202–2204.

906    Warkentine BE, Lavett Smith C, Rachlin JW. 1987. A Reevaluation of the Karyotype of the
907    Atlantic Silverside, *Menidia menidia*. Copeia. 1987:222-224.

908    Weisenfeld NI, Kumar V, Shah P, Church DM, Jaffe DB. 2017. Direct determination of diploid
909    genome sequences. Genome Res. 27:757–767.

910    Weissensteiner MH, Bunikis I, Catalán A, Francoijs K-J, Knief U, Heim W, Peona V, Pophaly S,
911    Sedlazeck FJ, Suh A et al. 2020. Discovery and population genomics of structural variation in a
912    songbird genus. Nat. Commun. 11:3403.

913    Wellenreuther M, Bernatchez L. 2018. Eco-Evolutionary Genomics of Chromosomal Inversions.
914    Trends Ecol. Evol. 33:427–440.

915    Wilder AP, Palumbi SR, Conover DO, Therkildsen NO. 2020. Footprints of local adaptation
916    span hundreds of linked genes in the Atlantic silverside genome. Evol Lett. 4:430–443.

917    Yuan Z, Liu S, Zhou T, Tian C, Bao L, Dunham R, Liu Z. 2018. Comparative genome analysis
918    of 52 fish species suggests differential associations of repetitive elements with their living
919    aquatic environments. BMC Genomics. 19:141.

920    Zhang D-C, Guo L, Guo H-Y, Zhu K-C, Li S-Q, Zhang Y, Zhang N, Liu B-S, Jiang S-G, Li J-T.
921    2019. Chromosome-level genome assembly of golden pompano (*Trachinotus ovatus*) in the
922    family Carangidae. Sci Data. 6:216.

923